

Jordi Alvarez, Núria Castell
 Departament de Llenguatges i Sistemes Informàtics,
 Universitat Politècnica de Catalunya

Pau Gargallo 5, 08028 Barcelona
 jalvarez, castell@lsi.upc.es

Resumen: en este artículo abordamos el proceso de validación de especificaciones preliminares escritas en lenguaje natural desde una nueva perspectiva: el uso de técnicas de aprendizaje automático. Describimos un sistema que a partir de ejemplos de especificaciones correctas construye o refina un modelo que nos permitirá validar nuevas especificaciones.

1. Introducción

La fase de ingeniería de requerimientos es una de las más importantes dentro del proceso de desarrollo del software. Se ha comprobado en diversos estudios que la mayoría de los errores detectados en la fase de codificación y mantenimiento provienen de una mala especificación [Basili:84, Endres:75]. Además, el coste de corregir un error de especificación se incrementa desmesuradamente cuando se detecta en las fases de codificación o mantenimiento [Pressman:92]. Esta fue una de las principales motivaciones por las que surgió el proyecto LESD [Castell:94]. Éste define cinco factores que nos permitirían validar la calidad de las especificaciones: trazabilidad, completitud, consistencia, verificabilidad y modificabilidad. LESD es continuado por el proyecto CICYT TIC93-0420 en el cual se enmarca este trabajo.

El artículo intenta dar una nueva perspectiva a la validación de especificaciones preliminares escritas en lenguaje natural. Tanto el hecho de que se trate de especificaciones preliminares como el que éstas estén escritas en lenguaje natural hace que sea realmente difícil dar una definición clara sobre los criterios que hacen que una especificación sea válida o no. Por ello proponemos usar Técnicas de aprendizaje automático para adquirir, a partir de ejemplos correctos de especificaciones, un modelo de especificación válida.

Para la representación de las especificaciones usamos una red semántica. Es un mecanismo que permita representar fácilmente estructuras jerárquicas y la mayoría de problemas que se presentan en ingeniería del software son de ese tipo [Maiden:95]; y además las redes semánticas, con la ayuda de los tres mecanismos básicos de representación (agregación, clasificación y especialización) de [Borgida:86] constituyen quizás el mecanismo de representación de conocimiento más adecuado para este tipo de problemas [Devanbu:92].

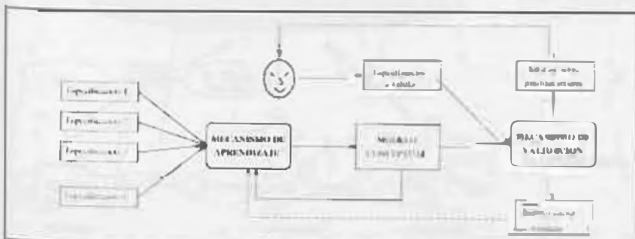


Figura 1: Esquema general del proceso de aprendizaje y validación

Uso de Técnicas de Aprendizaje para la Validación de Especificaciones

De este modo, a partir de especificaciones válidas, que se pueden ver como compuestas por una serie de descripciones de objetos, abstraeremos una serie de conceptos basándonos en esas descripciones. Clasificaremos cada uno de estos objetos como instancia de uno de los conceptos creados, y a partir de ellos aprenderemos una serie de propiedades que podemos interpretar como la descripción genérica de cada concepto. Posteriormente, para validar una especificación nos bastaría con clasificar los objetos cuya descripción constituye la especificación y comprobar si cumplen las propiedades correspondientes al concepto bajo el cual han sido clasificados.

En resumen, lo que pretendemos es generar o enriquecer una base de conocimiento a partir de especificaciones ya validadas; y posteriormente usar esta base de conocimiento para validar (y posiblemente también ayudar en la corrección de) otras especificaciones.

En la figura 1 podemos ver un esquema general del proceso. El aprendizaje puede ser para crear un modelo a partir de cero o bien para enriquecer un modelo que ya tuviéramos previamente. En cualquier caso se trata de un proceso automatizado en el que el ingeniero de software sólo se encarga de proporcionar los ejemplos. En cambio la validación no es un proceso totalmente automático: tras detectar los posibles errores, se genera un informe que pasa por el ingeniero de software, quien decide si se trata realmente de errores y cómo solucionarlos. Finalmente, una vez que la especificación está validada, puede ser utilizada por el mecanismo de aprendizaje para refinar el modelo conceptual que ya teníamos.

En este artículo nos ocupamos únicamente de los procesos de aprendizaje y validación una vez ya tenemos la representación conceptual de las especificaciones en nuestra red semántica. El paso de las especificaciones tal como las recibimos en lenguaje natural a dicha representación conceptual se describe en [Toussaint:92]; mientras que en [Castell:95] podemos encontrar una descripción global del proceso de control de especificaciones en el cual está inmerso el trabajo descrito aquí. En la sección 2 se comenta la importancia de las estructuras, el elemento básico a partir del cual realizamos el aprendizaje. En la sección 3 distinguimos los diferentes niveles de conocimiento que tenemos en una especificación de requerimientos y sobre los cuales podemos aprender. En la sección 4 profundizamos un poco los mecanismos de

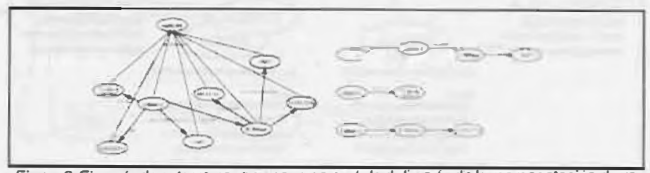


Figura 2: Ejemplo de estructuras básicas a partir de la definición de la representación de un libro y su autor en la red semántica

En la **figura 3** podemos observar las propiedades definitorias y asociadas de cuatro conceptos formados a partir de estructuras simples extraídas de la descripción de LIBRO-1.

3. Conocimiento y Metaconocimiento

Hasta ahora hemos hablado de agrupar en conceptos los objetos cuyas descripciones componen las especificaciones; pero no del tipo de conocimiento que el sistema debe aprender. Tipicamente, una especificación contiene dos niveles de información diferentes.

La información sobre objetos tangibles describe la realidad en sí (aquello que tradicionalmente en las redes semánticas se ha dado en llamar token y que en {Dardenne:93} constituye el nivel de instancias). P.ej. en una especificación de un sistema de gestión de una biblioteca, cuando hablamos de 'la biblioteca de la Facultad de Informática de Barcelona', estamos hablando de una entidad presente en la realidad del sistema especificado.

Las descripciones de conceptos del dominio (nivel de dominio en {Dardenne:93}). Siguiendo con el mismo ejemplo, en la especificación de un sistema de gestión de una biblioteca nos encontraremos seguramente con una descripción (adecuada a nuestro problema, eso sí) de librería, libro, etc.

A nivel de representación, el mecanismo de clasificación nos permite distinguir entre estos dos niveles de conocimiento e ir más allá. Si el sistema de representación es lo bastante flexible para tratar los conceptos como simples objetos, puede hacer que cualquier concepto, como objeto que es, sea a la vez instancia de otro concepto. Sistemas como TELOS {Mylopoulos:90} o KAOS {Dardenne:93, Lamsweerde:95} usan esto para colocar un nivel de metaconocimiento por encima del nivel de los conceptos pertenecientes al dominio de aplicación.

La información de este metanivel constituye un metamodelo que permite imponer ciertas restricciones sobre el nivel de conocimiento básico de igual forma que éste hace con el nivel de los tokens. Este metanivel constituye pues una herramienta potente para encontrar incompletitudes e inconsistencias en el nivel de descripción del dominio (el modelo).

El aprendizaje nos permite abstraer propiedades que cumplen los objetos de un nivel a partir de objetos del nivel inferior; mientras que el mecanismo de clasificación nos permite tener un modelo de un nivel inferior a partir de su inmediato superior. Debemos ser capaces de aprender tanto sobre los tokens

como sobre el modelo. Para ello, seguiremos la misma filosofía que los sistemas de representación mencionados. Tendremos tres niveles: los tokens, el conocimiento del dominio y el metanivel. En {Mylopoulos:90} se menciona la posibilidad de tener más de tres niveles de conocimiento, pero los realmente útiles al menos en este caso son los tres primeros, pues se trata de modelizar la información que aparece en la especificación y dicha información pertenece a los dos primeros niveles de tokens y de conceptos básicos (la modelización de un nivel viene dada por el nivel inmediatamente superior, el único que puede imponer condiciones sobre él: basta añadir un tercer nivel, que no necesitamos modelizar sino simplemente utilizar).

Así podremos separar el proceso global de aprendizaje en dos niveles distintos: el aprendizaje sobre el **conocimiento del dominio** (a partir del nivel de tokens) y el aprendizaje del **metamodelo** (realizado a partir del conocimiento del dominio). Según nos interese, podemos aplicar el primero, con lo que completaremos y aprenderemos las propiedades de la jerarquía de conceptos del dominio; aplicar el segundo, con lo que generaremos un metamodelo (o completaremos el que ya teníamos); o bien aplicar los dos. De esta forma, tal como se muestra en la **figura 4**, se produce un flujo de información bidireccional (por un lado el aprendizaje y por otro un supuesto sistema de tipos regido por el mecanismo de clasificación) entre el nivel de las instancias y el nivel del dominio {Feather:93}; y entre este último y el metanivel.

4. Generación de Conceptos

La base sobre la cual vamos a construir nuestra jerarquía (bien sea la correspondiente al modelo o la del metamodelo) son las propiedades estructurales de los objetos. Para ello, los conceptos que crearemos se definirán a partir de cadenas estructurales conceptualizadas; entendiendo por cadena estructural una serie de nodos ordenados y unidos cada uno con el siguiente (excepto el último) mediante una relación; y por **conceptualizadas** que aparezca en la cadena, no el objeto que aparece en la especificación inmerso en la estructura, sino el concepto del cual es instancia. En la **figura 5** podemos ver las estructuras simples de la **figura 2** conceptualizadas.

El hacer el paso de la conceptualización se debe a que el conocimiento que queremos aprender (recordemos que queremos aprender un modelo para una serie de objetos) está a nivel de conceptos y no a nivel de instancias. Por lo tanto sería una pérdida de tiempo intentar razonar sobre las instancias cuando aportan el mismo conocimiento que los conceptos correspon-

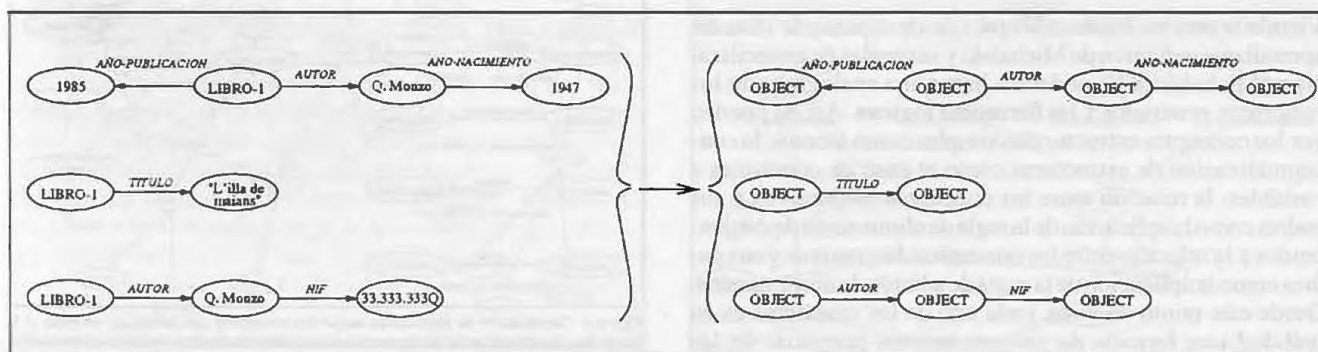


Figura 5: Conceptualización de las estructuras de la fig. 2 teniendo en cuenta que todos los objetos son instancia de objeto

dientes. Son los conceptos quienes dan un significado a sus instancias, que por sí mismas no lo tienen (al principio del proceso de aprendizaje, si partimos de cero, la jerarquía tendría un único concepto al que podemos dar el nombre genérico de object, o meta-object según el caso. Eso significa que al principio medimos todos los objetos por el mismo rasero, y object aparecería como único nodo en todas las posiciones de todas las cadenas estructurales. Esta situación es acorde con que lo único que tenemos para aprender es la estructura si en un principio no tenemos ninguna jerarquía).

La conceptualización de las estructuras es lo que nos permite sacar provecho de los mecanismos de clasificación y generalización. Respecto al primero, su uso en la conceptualización de estructuras es evidente al cambiar los objetos por los conceptos de los que son instancia. En cuanto al segundo, la posterior clasificación en la jerarquía del concepto generado hace que obtengamos estructuras más generales que otras simplemente porque los conceptos y relaciones que participan en unas son más generales que los que participan en las otras. Así por ejemplo, en la **figura 6** podemos ver como los conceptos STRC-1 Y STRC-3 son más generales que STRC-1 Y STRC-3 porque repositorio es más general que biblioteca.

4.1. Tipos de Conceptos

Para construir una jerarquía a partir de cadenas estructurales tenemos tres tipos de conceptos distintos: **conceptos estructurales simples, conceptos conjuntivos y conceptos disyuntivos**. Los primeros vienen definidos por una única propiedad estructural que es una cadena estructural conceptualizada. Los conceptos conjuntivos, como su nombre indica, vienen definidos por la conjunción de las propiedades estructurales que definen a los padres. Y por último, los conceptos disyuntivos se definen a partir de la disyunción de las propiedades estructurales que definen a sus hijos.

Los conceptos estructurales simples constituyen las piezas básicas sobre las que los conceptos conjuntivos y disyuntivos nos permitirían hacer especializaciones y generalizaciones. Con estos tres tipos de conceptos combinados adecuadamente podremos obtener fácilmente conceptos con definiciones estructurales complejas. En la **figura 7** podemos ver cómo sirviéndonos de conceptos conjuntivos y estructurales simples construimos una jerarquía en la que especificamos las descripciones mínimas de distintos tipos de publicaciones.

4.2. Analogía con el Aprendizaje Inductivo

Viendo la generación de conceptos desde el punto de vista del aprendizaje inductivo de Michalski y sus reglas de generalización [Michalski:83], podemos hacer una analogía entre los conceptos generados y las **fórmulas lógicas**. Así, se pueden ver los conceptos estructurales simples como átomos, la conceptualización de estructuras como el paso de constantes a variables, la relación entre los conceptos conjuntivos y sus padres como la aplicación de la regla de eliminación de conjuntivos y la relación entre los conceptos disyuntivos y sus padres como la aplicación de la regla de adición de disyuntivos. Desde este punto de vista, cada uno de los conceptos es en realidad una función de reconocimiento potencial de las propiedades que queremos aprender. No nos debe importar

que algunas reglas de generalización de las usadas no preserven la equivalencia lógica entre una fórmula y su generalización (de las tres mencionadas, la única que no la preserva es la de paso de constantes a variables) [Kodratoff:88]. El aprendizaje real de las propiedades se realiza sobre los conceptos generados en una fase posterior; y puede revelarnos algunos de estos conceptos como totalmente inútiles. Lo que sí nos debe preocupar es que hayamos generado como mínimo todos aquellos conceptos que representen verdaderamente funciones de reconocimiento sobre las propiedades que queremos aprender.

Como veremos, el mecanismo de aprendizaje de las propiedades requiere que el nodo relativo al concepto se encuentre explícitamente representado en la red semántica. De no ser así, el aprendizaje se realizaría sobre conceptos representativos de funciones de reconocimiento que por muy parecidas que sean a la verdadera, caerían en la incompletitud o inconsistencia.

La única manera que tenemos de asegurar que todos los conceptos que representen funciones de reconocimiento que nos interesan sean generados es creando todos los conceptos posibles. Pero esto es computacionalmente imposible dado el elevado número de estructuras distintas con que nos podemos encontrar. Por ello, no nos queda más remedio que utilizar heurísticas que nos guíen en la generación de conceptos.

4.3. Modelo Heurístico de Generación de Conceptos

Este es un punto que todavía tenemos que explorar en profundidad: si bien hasta este momento hemos utilizado un sistema bastante sencillo que funciona aceptablemente para conocimiento generado aleatoriamente a partir de un metamodelo previo que es destruido justo antes de empezar el aprendizaje (parecido al utilizado en KAOS [Dardenne:93]).

El proceso consiste en generar tantas estructuras simples como aparezcan en las descripciones de ejemplos que tengamos. Para no disparar la complejidad de los cálculos las estructuras generadas no deben sobrepasar una longitud máxima preestablecida. De todas maneras, si con estructuras de una determinada longitud no tuviéramos suficiente, la

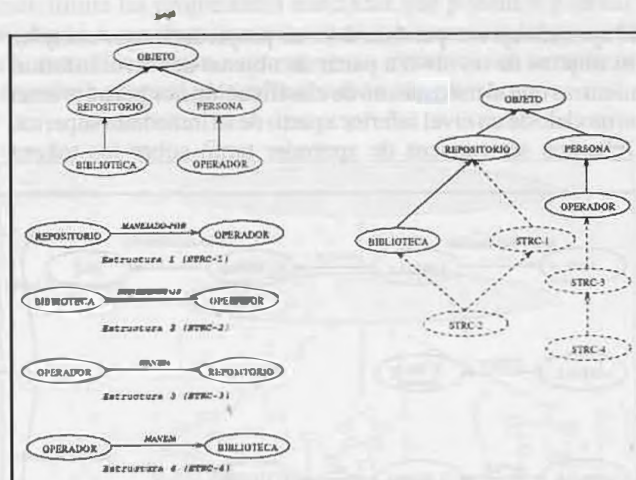


Figura 6: Clasificación de estructuras según los conceptos que participan en ellas. A la izquierda, una simplificación de la jerarquía inicial y las estructuras que vamos a introducir. A la derecha podemos ver como queda la jerarquía después de introducir los nuevos conceptos y clasificarlos automáticamente.

iteración del proceso nos permitiría, tal como veremos más adelante, acceder a estructuras más largas de forma automática.

Tras obtener las estructuras simples, usamos conceptos conjuntivos para generar una estructura compleja para cada descripción de objeto que tengamos y luego clasificamos cada objeto como instancia del concepto correspondiente. La herencia múltiple puede realizarse siempre a nivel conceptual; lo que permite al mecanismo de aprendizaje abstraer propiedades relacionadas con la conjunción de varios conceptos simples.

Para terminar de construir la jerarquía usamos *técnicas de clustering* {Jain:88} sobre las estructuras básicas generadas. La distancia de un concepto a un *cluster* se calcula como la máxima de sus distancias a todos los elementos del cluster. La distancia entre distintos conceptos se calcula a partir de la descripción de cada concepto, que viene dada por el conjunto de propiedades asociadas a él. Por lo tanto, es necesario haber realizado el aprendizaje de propiedades, explicado en la próxima sección, sobre los conceptos a los que se aplica el clustering.

De forma intuitiva, la distancia (y también la similitud) entre las diferentes descripciones de conceptos se calcula a partir de la información común (propiedades asociadas iguales) que tienen, penalizando la ausencia de información. Esta ausencia de información se mide a partir del total de tipos de propiedades asociadas que puede tener un concepto. Su penalización hace que la supuesta distancia entre dos elementos no sea tal distancia sino simplemente un heurístico; ya que una de las propiedades que debe cumplir es que la distancia entre un elemento y él mismo sea siempre nula {Kodratoff:88}. En nuestro caso sólo sería nula si su descripción contiene todos los tipos de propiedades asociadas existentes. De todas maneras, esta penalización proporciona una medida heurística que se muestra como la mejor de las que hemos probado. Además, fuerza a construir una jerarquía de forma que los objetos con mucha información (ya sean conceptos o instancias de ellos) estarían situados en la parte profunda de la jerarquía, mientras que los que tienen poca información estarían situados más arriba.

No generamos un único nivel de clusters sino que empezamos con un umbral bajo de distancia (podemos empezar con cero, por ejemplo) y vamos iterando el proceso mientras

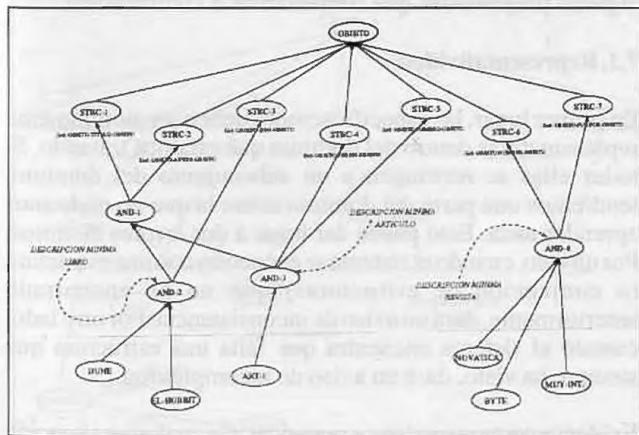


Figura 7: Distinción de distintos tipos de publicación mediante el uso de propiedades estructurales definitorias

incrementamos este umbral hasta llegar a 1 (la máxima distancia entre dos conceptos). De esta forma, acabaremos generando un único cluster en la parte más alta de la jerarquía que se correspondería con la raíz de ésta, con lo cual la habremos completado.

5. Aprendizaje de las Propiedades

Una vez que ya tenemos los nodos explícitamente representados en la red semántica, procedemos con el aprendizaje de las propiedades estructurales asociadas a cada concepto. Esta fase se puede realizar bien cuando los conceptos están ya todos generados, o bien incrementalmente a medida que se van generando (de hecho, las técnicas de clustering del apartado anterior requieren que en los conceptos sobre los que se aplica el clustering ya se haya hecho la abstracción de las propiedades asociadas).

Como se puede observar en la figura 8, las propiedades, tanto definitorias como asociadas y por complejas que sean, se pueden reducir a propiedades asociadas más simples entre conceptos más complejos. Por ello podemos reducir la representación de tales propiedades a la representación de simples *links* o relaciones (con un peso que indica el grado en que la propiedad se cumple). Y por tanto, podemos reducir también el aprendizaje al aprendizaje de simples *links*. En la figura se puede observar como tras la primera iteración del proceso (dibujada con trazo de línea continua) se genera la propiedad definitoria OBJETO-AUTOR-OBJETO para el concepto STRC-1. En la segunda iteración (trazo punteado) generamos un nuevo concepto que especializa la propiedad definitoria de STRC-1 a STRC-1-AUTOR-AND-1. Y como AND-1 es un concepto definido a partir de propiedades estructurales, la propiedad definitoria de STRC-5 equivale a la conjunción de: STRC-1-AUTOR-OBJETO-NIF-OBJETO, STRC-1-AUTOR-OBJETO-TELEFONO-OBJETO y STRC-1-AUTOR-OBJETO-NOMBRE-OBJETO. Lo mismo ocurre para las propiedades asociadas ya que aparecen en ellas cada vez nodos con definiciones estructurales más complejas a base de iterar el proceso de aprendizaje.

En definitiva, el aprendizaje consiste pues en decidir en qué grado los conceptos generados están en relación con otros conceptos. Para ello, inspirándonos libremente en el modelo propuesto por Shastri en {Shastri:88}, damos a cada *link* entre dos conceptos el valor correspondiente al porcentaje de instancias del primer concepto que están relacionadas mediante

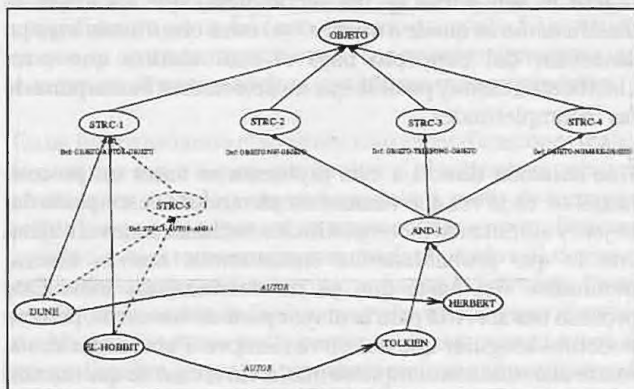


Figura 8: Iteración del proceso de generación de conceptos

este tipo de *link* a instancias del segundo concepto. De esta forma, suponiendo que los nodos de nuestra red semántica se activaran y que cuando se activara un objeto se activaran todos los conceptos del cual es instancia, el peso de un *link* es una medida de la frecuencia de activación del nodo destino se activa cuando el nodo origen está activado; lo cual se ajusta extremadamente bien a una interpretación Hebbiana de los pesos sinápticos en las redes neuronales [Shastri:88]. Pero a diferencia de Shastri, nuestros *links* están etiquetados (sus etiquetas son las relaciones); por lo que podemos tener varios *links* diferentes entre los mismos dos conceptos.

6. Validación de Especificaciones

Una vez tenemos realizado el proceso de aprendizaje, podemos empezar a validar especificaciones. La validación es sencilla. Si interpretamos la especificación que vamos a validar como un conjunto de descripciones de los objetos que intervienen en el sistema especificado, sólo necesitamos clasificar en la jerarquía aprendida cada uno de estos objetos y posteriormente comparar sus propiedades con las del concepto bajo el cual ha sido clasificado.

Las propiedades definitorias de los conceptos servirían para clasificar los objetos, mientras que las propiedades asociadas servirían para verificar que su descripción sea completa y consistente.

Si un concepto tiene alguna propiedad estructural asociada que alguna de sus instancias no tiene, habremos detectado una posible incompletitud (mas adelante ahondaremos en el porqué del término posible). La razón es sencilla: si un grupo de objetos que tienen en común parte de su estructura (las propiedades estructurales definitorias del concepto que los agrupa) tienen todos ellos (o la mayoría) cierta propiedad y nos encontramos con otro objeto que tiene en común toda la estructura menos esta última propiedad, sospecharemos la posibilidad de que la descripción de este objeto sea incompleta. La evidencia que tengamos de esto dependería del porcentaje de instancias que cumplan la propiedad. Recíprocamente, si un objeto tiene alguna propiedad que el concepto bajo el cual ha sido clasificado no tiene, sospecharemos la existencia de una posible inconsistencia o bien que la descripción del objeto contiene información redundante.

Pero no todo termina ahí. Hay otro caso que permite detectar una posible incompletitud. ¿Qué pasa si la descripción de un objeto es tan pobre (y tan incompleta) que hace que la clasificación se quede a medias? Al estar clasificado bajo un antecesor del concepto bajo el cual tendría que estar clasificado, es muy posible que no detectemos buena parte de las incompletitudes.

Una solución directa a este problema es hacer un proceso iterativo: cada vez que vamos completando la descripción del objeto y eliminando incompletitudes, reclasificamos el objeto, con lo que probablemente encontremos nuevas incompletitudes. Así hasta que ya no encontramos más. Este proceso nos serviría para la mayor parte de los casos, pero no podemos asegurar que nos sirva siempre. Para asegurarnos, existe una solución complementaria. En el caso de que nuestro objeto haya sido clasificado bajo un concepto que no tiene

instancias directas, sino que sus instancias son instancias de algún concepto más específico (lo que en terminología de lenguajes orientados a objetos se llamaría una clase abstracta), podemos deducir que su descripción es tan incompleta que no permite una clasificación total del objeto. De esta forma, si el llamémosle-concepto abstracto no permite iniciar el proceso iterativo que acabamos de comentar, podemos intentar guiarnos por las propiedades definitorias de los conceptos hijos. Es decir; intuimos que el objeto no puede ser instancia directa del concepto supuestamente abstracto, por lo que debe ser instancia (directa o no, da igual) de alguno de sus hijos. Por lo tanto, intentamos averiguar qué propiedad o propiedades definitorias de los hijos debería cumplir el objeto y evidentemente no cumple (si no, la clasificación hubiera continuado).

En la figura 9 podemos ver como una descripción incompleta de un libro se clasifica bajo AND-1. Al no tener éste ninguna instancia directa, el sistema sospecharía que la descripción es incompleta y propondría al usuario completar la descripción de EL-HOBBIT con propiedades correspondientes a AND-2 y AND-3 (según lo que conoce el sistema, la descripción del objeto puede evolucionar hacia un libro o hacia una revista).

Mientras en los otros casos de detección de posibles incompletitudes e inconsistencias, estaba claro cuales eran las propiedades estructurales en cuestión, aquí tan sólo sabemos que se trata de un subconjunto de las propiedades estructurales que diferencian a los hijos del concepto abstracto de éste. Es el ingeniero de software quien debe decidir cuales son las propiedades que faltan (si es que realmente falta alguna).

7. Conjunto de Entrenamiento

En cualquier método de aprendizaje a partir de ejemplos, tiene vital importancia la elección del conjunto de ejemplos a partir de los cuales vamos a realizar el entrenamiento (conjunto de entrenamiento). En nuestro caso, el conjunto de entrenamiento está formado por especificaciones. Después de hablar de la validación en el apartado anterior, analizaremos qué propiedades deben cumplir las especificaciones del conjunto de entrenamiento para que el aprendizaje y la posterior validación sean eficaces. Y para que el aprendizaje sea eficaz, todas en conjunto y cada una en particular deberán cumplir algunas propiedades que comentamos a continuación.

7.1. Representatividad

En primer lugar, las especificaciones deben ser ampliamente representativas dentro del dominio que estamos tratando. Si todas ellas se restringen a un subconjunto del dominio, tendremos una parte del dominio sobre la que no podremos aprender nada. Esto puede dar lugar a dos errores distintos. Por un lado, cuando el sistema se encuentra con una estructura (o conjunción de estructuras) que no ha encontrado anteriormente, dará un aviso de inconsistencia. Por otro lado, cuando el sistema encuentra que falta una estructura que siempre ha visto, dará un aviso de incompletitud.

Evidentemente no podemos pretender que cualquier situación posible en nuestro dominio esté presente en alguna de las

especificaciones de ejemplo, con lo que no podemos asegurar que las incompletitudes e inconsistencias detectadas por el sistema lo sean de verdad. Por ello, la decisión final la debe tomar siempre el ingeniero de software {Reubenstein:91} y debemos entender los avisos del sistema de validación como lo que son: avisos sobre posibles incompletitudes o inconsistencias que el sistema deja decidir si realmente lo son o no al buen juicio del ingeniero de software.

Todo esto no quiere decir que el sistema presentado no sirva para nada. Después de un extenso análisis de problemas propios de ingeniería del software, {Maiden:95} argumenta que la mayor parte de ellos pertenecen a un conjunto tratable de clases jerárquicas. Si las abstracciones realizadas en la fase de aprendizaje son suficientemente buenas, y para un mismo dominio de aplicación y entorno de desarrollo, el número de situaciones nuevas que puede presentar una especificación respecto a un conjunto de ellas (el de entrenamiento) decrece rápidamente si el tamaño del conjunto crece. Además, una vez validadas, las nuevas especificaciones se pueden ir incorporando al conjunto de entrenamiento, con lo que los posibles avisos erróneos se van corrigiendo y la fiabilidad del sistema de validación va en aumento.

7.2. Consistencia

Evidentemente, cada una de las especificaciones debe ser consistente. De no ser así, el sistema aprenderá las inconsistencias de los ejemplos y el resultado de la validación perderá cualquier valor que pudiera tener.

7.3. Minimalidad

Importa que las especificaciones que formen parte del conjunto de entrenamiento sean mínimas. Así nos aseguramos que todas las propiedades que aprende el sistema son estrictamente necesarias en la especificación. Esto es especialmente importante para comprobación de la completitud, ya que justifica que el sistema pueda asegurar que existe una incompletitud cuando detecta la ausencia de una de esas propiedades. De no ser mínimas las especificaciones, podríamos aprender propiedades innecesarias, lo que haría que el sistema pudiera detectar incompletitudes allí donde no las hay.

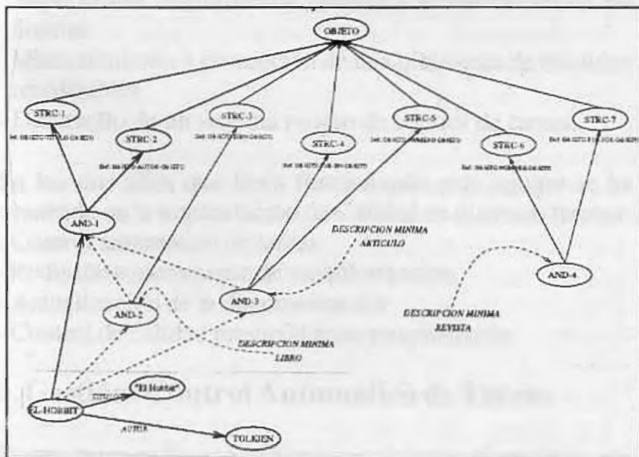


Figura 9. Validación de una descripción de libro a la cual le falta el atributo ISBN

7.4. Completitud

Recíprocamente, también es importante que las especificaciones sean completas. Si no, el sistema podría aprender descripciones incompletas de objetos y por tanto detectar menos incompletitudes de las realmente existentes. Además, la incompletitud de las especificaciones del conjunto de entrenamiento también afectaría a la detección de inconsistencias. El sistema no aprendería algunas de las propiedades necesarias, y cuando se las encontrara, las reportaría como inconsistencias.

7.5. Relajación de las Condiciones

Como para realizar el aprendizaje de las propiedades asociadas a cada concepto no se utilizan técnicas de aprendizaje simbólico sino que únicamente se mantiene un ratio del porcentaje de instancias que cumplen dichas propiedades, el sistema es bastante robusto y resistente a errores. Por ello, no es necesario que las especificaciones del conjunto de entrenamiento cumplan al cien por cien las condiciones especificadas. El sistema construirá un modelo de validación más o menos correcto aunque las especificaciones no sean totalmente correctas. Eso sí, cuantos más errores, más ejemplos necesitaremos para obtener una validación fiable.

8. Conclusiones y trabajo Futuro

Hemos presentado en este artículo un sistema que a partir de especificaciones ejemplares es capaz de abstraer un modelo de especificación para comprobar la completitud y la consistencia de nuevas especificaciones.

La construcción manual de una base de conocimiento que permita validar especificaciones resulta altamente inviable y costosa debido a la complejidad del problema (sobre todo en lo que se refiere a la completitud), a la total vinculación de los términos completitud y consistencia al dominio de aplicación y a la metodología usada para el desarrollo del software en cada caso concreto, así como al elevado grado de informalidad que implica el que se trate de especificaciones preliminares escritas en lenguaje natural. Por ello, el uso de Técnicas de aprendizaje automático para la construcción de dicha base de conocimiento parece una buena solución.

Por el momento hemos realizado pruebas satisfactorias, generando automáticamente conocimiento a partir de un metamodelo parecido al usado en el sistema KAOS y aprendiendo a partir del conocimiento generado. El resultado es que se aprenden (entre otros) los mismos conceptos que estaban en el metamodelo y con las mismas propiedades.

En un futuro próximo planeamos usar especificaciones reales, lo cual nos serviría para comprobar si el sistema es capaz de abstraer un modelo y un metamodelo a partir de ejemplos sobre los que no tenemos un metamodelo previo. Esto nos ayudará a hacer un estudio más profundo sobre el conjunto de heurísticas que utilizamos para generar la jerarquía. Puede resultar interesante añadir algunas reglas más de generalización de Michalski, como por ejemplo la de subir en el árbol de generalización, lo que nos permitiría realizar automáticamente la generalización de la figura 6.

Por otro lado, para incrementar la potencia del sistema de representación, puede ser interesante introducir en la propia estructura conceptos reflexivos sobre ésta misma. Por ejemplo, podríamos introducir el concepto SELF para referirse de una manera estándar al primer elemento de la estructura. Esto permitiría entre otras cosas descubrir propiedades de las relaciones temporales, comprobando que la cadena estructural CONCEPTO—ANTES—CONCEPTO—ANTES—SELF nunca se cumple y por tanto es inconsistente que A pase antes que B si B pasa antes que A. A un nivel más abstracto podríamos aprender qué relaciones son inversas de otras mediante la cadena: CONCEPTO—RELACION—CONCEPTO—INVERSA—SELF.

Por último, también pensamos trabajar sobre un sistema que permita ir eliminando todo el conocimiento generado que no nos sea útil. El sistema genera gran cantidad de conceptos y conexiones que después resultan inútiles. Si para problemas reales no se elimina buena parte de esta 'basura', la cantidad de espacio ocupado puede ser desmesurada.

Bibliografía

- V. Basili, B. Perricone.** *Software errors and complexity: An empirical investigation.* Communications of the ACM, 27(1):42—52, Jan. 1984.
- A. Borgida.** *Conceptual modeling of information systems.* In M. Brodie and J. Mylopoulos, editors, Knowledge Base Management Systems, chapter 31, pages 461—469. Springer-Verlag, 1986.
- N. Castell, A. Hernández.** *Filtering software specifications written in natural language.* In Proceedings of 7th Portuguese Conference on Artificial Intelligence, (EPIA'95), pages 447—455, Madeira, Portugal, Oct. 1995. LNAI-990, Springer Verlag.
- N. Castell, O. Slavkova, Y. Toussaint, A. Tuells.** *Quality control of software specifications written in natural language.* In Proceedings of 7th International Conference on Industrial & Engineering Applications of AI & Expert Systems (IEA-AIE'94), pages 37—44, Austin, Texas, June 1994. Gordon and Breach Science Publishers.
- A. Dardenne, A. van Lamsweerde, S. Fickas.** *Goal-directed requirements acquisition.* Science of Computer Programming, 20:3—50, 1993.
- P. Devanbu, R. Brachman, P. Selfridge, B. Ballard.** *Lassie: A knowledge-based software information system.* Communications of the ACM, 34(5):35—49, May 1991.
- A. Endres.** *An analysis of errors and their causes in system programs.* IEEE Transactions on Software Engineering, 1(2):140—149, June 1975.
- M. Feather.** *Requirements reconnoitering at the juncture of domain and instance.* In Proceedings of 1st International Symposium on Requirements Engineering, pages 73—76, San Diego, California, Jan. 1993. IEEE Computer Society Press.
- A. Jain, R. Dubes.** *Algorithms for Clustering Data.* Prentice-Hall, 1988.
- Y. Kodratoff.** *Introduction to Machine Learning.* Morgan Kaufman Publishers, 1988.
- N. Maiden, P. Mistry, A. Sutcliffe.** *How people categorise requirements for reuse: a natural approach.* In Proceedings of 2nd International Symposium on Requirements Engineering, pages 194—203, York, England, Mar. 1995. IEEE Computer Society Press.
- R. Michalski.** *A theory and methodology of inductive learning.* In R. Michalski, J. Carbonell and T. Mitchell, editors, Machine Learning, An Artificial Intelligence Approach, pages 83—130. Morgan Kaufman Publishers, 1983.
- J. Mylopoulos, M. Jarke, A. Borgida, M. Koubarakis.** *Telos: Representing knowledge about information systems.* ACM Transactions on Information Systems, 8(4):325—362, Oct. 1990.
- R. Pressman.** *Software Engineering: A Practitioner's Approach.* McGraw-Hill, New York, NY, USA, third edition, 1992.
- H. Reubenstein, R. Waters.** *The requirements apprentice: Automated assistance for requirements acquisition.* IEEE Transactions on Software Engineering, 17(3):226—240, Mar. 1991.
- L. Shastri.** *A connectionist approach to knowledge representation and limited inference.* Cognitive Science, 12:331—392, 1988.
- G. Spanoudakis, P. Constantopoulos.** *Measuring similarity between software artifacts.* In Proceedings of 6th International Conference on Software Engineering and Knowledge Engineering, Jurnala, Latvia, June 1994.
- Y. Toussaint.** *Méthodes Informatiques et Linguistiques pour l'aide à la Spécification de Logiciel.* PhD thesis, Paul Sabatier University, Toulouse, 1992.
- A. van Lamsweerde, R. Darimon, P. Massonet.** *Goal-directed elaboration of requirements for a meeting scheduler: Problems and lessons learnt.* In Proceedings of 2nd International Symposium on Requirements Engineering, pages 194—203, York, England, Mar. 1995. IEEE Computer Society Press.

Agradecimientos

Este trabajo está parcialmente subvencionado por la CICYT (TIC93-420) y por la CIRIT (GRQ93-3015) y una beca predoctoral a Jordi Alvarez.